

Bigram Anchor Words Topic Model

Ashuha Arseniy¹ and Loukachevitch Natalia²

¹ Moscow Institute of Physics and Technology,
`ars.ashuha@pystech.edu`,

² Research Computing Center of Lomonosov Moscow State University,
`louk_nat@mail.ru`

Abstract. A probabilistic topic model is a modern statistical tool for document collection analysis that allows extracting a number of topics in the collection and describes each document as a discrete probability distribution over topics. Classical approaches to statistical topic modeling can be quite effective in various tasks, but the generated topics may be too similar to each other or poorly interpretable. We supposed that it is possible to improve the interpretability and differentiation of topics by using linguistic information such as collocations while building the topic model. In this paper we offer an approach to accounting bigrams (two-word phrases) for the construction of Anchor Words Topic Model.

Keywords: topic model, anchor words, bigram

1 Introduction

A probabilistic topic model is a modern statistical tool for document collection analysis that allows identifying a set of topics in the collection and describes each document as a discrete probability distribution over topics. The topic is meant a discrete probability distribution over words, considered as a thematically related set of words. Topic models are actively used for various applications such as text analysis [1,2,3], users' analysis [4], information retrieval [5,6].

To recognize hidden topics, standard algorithms of topic modeling such as PLSA or LDA [7,8], take into account only the frequencies of words and do not consider the syntactic structure of sentences, the word order, or grammatical characteristics of words. Neglect of the linguistic information causes the low interpretability and the low degree of topic differentiation [6], which may hinder the use of topic models. If a topic has the low interpretability then it may seem as a set of unrelated words or a mixture of several topics. It is difficult to differentiate topics when they are very similar to each other.

One of the approaches that improves the interpretability of the topics is proposed in [9,10] and is called Anchor Words. This approach is based on the assumption that in each topic there exists a unique word that describes the topic, but this approach is also built on word frequencies.

In this paper we put forward a modification of the Anchor algorithm, which allows us to take into account collocations when building a topic model. The ex-

periments were conducted on various text collections (Banks Articles, 20 Newsgroups, NIPS) and confirmed that the proposed method improved the interpretability and the uniqueness of topics without downgrading other quality measures.

The paper is organized as follows. Section 2 reviews similar work. Section 3 describes the metrics used for evaluating the quality of topic models. In Section 4, we propose a method that allows us to take into account collocations in the Anchor Words topic model.

2 Related work

2.1 Notation and basic assumptions

Many variants of topic modeling algorithms have been proposed so far. Researchers usually suppose that a topic is a set of words that describe a subject or an event; a document is a set of topics that have generated it. A **Topic** is a discrete probability distribution over words: topic $t = \{P(w|t) : w \in W\}$ [7,8]. In this notation, each word in each topic has a certain probability, which may be equal to zero. Probabilities of words in topics are usually stored in the matrix $\Phi_{W \times T}$. A **Document** is a discrete probability distribution over topics $P(t|d)$ [7,8]. These probabilities are represented as a matrix $\Theta_{T \times D}$.

In topic modeling, the following hypotheses are usually presupposed: a **Bag of words hypothesis** is the assumption that it is possible to determine which topics have generated the document without taking into account the order of words in the document; **Hypothesis of conditional independence** is the assumption that the topic does not depend on the document, the topic is represented by the same discrete distribution in each document which contains this topic. Formally, the probability of a word in the topic is not dependent on the document – $P(w|d, t) = P(w|t)$ [6,7,8]; **Hypothesis about the thematic structure of the document** assumes that the probability of a word in a document depends on the hidden topics that have generated the document, as, for example, in the simplest topic model:

$$p(w|d) = \sum_{t \in T} P(w|d, t)P(t|d) = \sum_{t \in T} P(w|t)P(t|d) \quad (1)$$

2.2 Specific Topic models

In this section we consider several well-known approaches to topic modeling.

Probabilistic Latent Semantic Analysis, PLSA was proposed by Thomas Hoffman in [7]. To build the model, he supposed to optimize the log-likelihood

with the restrictions of normalization and non-negativeness:

$$\log L(D, \Phi, \Theta) = \log \prod_{d \in D} \prod_{w \in d} p(w|d) \rightarrow \max_{\Phi, \Theta} \quad (2)$$

$$\phi_{wt} \geq 0; \sum_{w \in W} \phi_{wt} = 1; \theta_{td} \geq 0; \sum_{t \in T} \phi_{td} = 1 \quad (3)$$

To solve the optimization problem, it was proposed to apply the EM-algorithm, which is usually used to find the maximum likelihood estimate of probability model parameters when the model depends on hidden variables.

Latent Dirichlet allocation, LDA was proposed by David Blei in [8]. This paper introduces the generative model that assumes that the vectors of topics and the vectors of documents are generated from the Dirichlet distribution. For training the model, it was proposed to optimize the following function:

$$\log L(D, \Phi, \Theta) \prod_d \text{Dir}(\theta_d | \beta) \prod_t \text{Dir}(\phi_t | \alpha) \rightarrow \max_{\Phi, \Theta} \quad (4)$$

$$\phi_{wt} \geq 0; \sum_{w \in W} \phi_{wt} = 1; \theta_{td} \geq 0; \sum_{t \in T} \phi_{td} = 1 \quad (5)$$

To solve the optimization problem, the authors use the Bayesian inference, which leads to EM-algorithm similar to PLSA. Because of the factored-conditional conjugate prior distribution and the likelihood, the formula for the parameters update can be written explicitly.

Additive Regularization Topic Model was proposed by Konstantin Vorontsov in [6]. The "Additive Regularization Topic Model" generalizes LDA (the LDA approach can be expressed in terms of an additive regularization) and allows applying a combination of regularizers to topic modeling by optimizing the following functional:

$$\log L(\Phi, \Theta) + \sum_{i=1}^n \tau_i R_i(\Phi, \Theta) \rightarrow \max_{\Phi, \Theta} \quad (6)$$

$$\phi_{wt} \geq 0; \sum_{w \in W} \phi_{wt} = 1; \theta_{td} \geq 0; \sum_{t \in T} \phi_{td} = 1 \quad (7)$$

where, τ_i – weight of regularizer, R_i – regularizer.

To introduce regularizers, the Bayesian inference is not used. On the one hand, it simplifies the process of entering regularizers because it does not require the technique of Bayesian reasoning, on the other hand, the introduction of a new regularizer is an art, which is hard to formalize.

The paper [6] shows that the use of the additive regularization allows simulating reasonable assumptions about the structure of topics, which helps to improve some properties of a topic model such as interpretability and sparseness.

Anchor Words Topic Model was proposed by Sanjeev Arora in [9,10]. The basic idea of this method is the assumption that for each topic t_i there is an anchor word that has a nonzero probability only in the topic t_i . If one have the *anchor words* one can recover a topic model without EM algorithm.

The algorithm 1 consists of two steps: the search of anchor words and recovery of a topic model with anchor words. Both procedures use the matrix $Q_{W \times W}$ that contains joint probabilities of co-occurrence of word pairs $p(w_i, w_j)$, $\sum Q_{ij} = 1$. Let us denote row-normalized matrix Q as \hat{Q} , the matrix \hat{Q} can be interpreted as $\hat{Q}_{i,j} = p(w_j|w_i)$.

Algorithm 1 High Level Anchor Words

Input: collection D, number of topics $|T|$

Output: matrix Φ ;

- 1: $Q = \text{Word Co-occurences}(D)$
 - 2: $\hat{Q} = \text{Rows_normalized}(Q)$
 - 3: $\hat{Q} = \text{Random_projection}(\hat{Q})$
 - 4: $S = \text{FindAnchorWords}(\hat{Q}, |T|)$
 - 5: $\Phi = \text{RecoverWordTopic}(\hat{Q}, S)$
 - 6: **return** Φ
-

Let us denote indexes of anchor words $S = \{s_1, \dots, s_T\}$. The rows indexed by elements of S are special in that every other row of \hat{Q} lies in the convex hull of the rows indexed by the anchor words [10]. At the next step optimization problems are solved. It's done to recover the expansion coefficients of $C_{it} = p(t|w_i)$, and then using the Bayes rule we restore matrix $(p(w|t))_{W \times T}$. The search of anchor words is equal to the search for almost convex hull in the vectors of the matrix \hat{Q} [10]. The combinatorial algorithm that solves the problem of finding the anchor words is given in Algorithm 2.

Algorithm 2 The combinatorial algorithm FastAnchorWords

Input: dots $V = v_1, \dots, v_n$, dim of convex hull K , parameter of error ϵ ;

Output: $\{v'_1, \dots, v'_k\}$ – set of points which constitute the convex hull;

- 1: put v_i into random subspace V , $\dim V = 4 \log V / \epsilon^2$
 - 2: $S = \{s_0\}$, s_0 – point that has the largest distance to origin.
 - 3: **for all** i **do** 1 to K :
 - 4: denote point $\in V$ that has the largest distance to $\text{span}(S)$ as s_i
 - 5: $S = S \cup \{s_i\}$
 - 6: **for all** i **do** 1 to K :
 - 7: denote point $\in V$ that has the largest distance to $\text{span}(S \setminus \{s_i\})$ as s'_i
 - 8: update s_i on s'_i
 - return** S
-

2.3 Integration of n-grams into topic models

The above-discussed topic models are based on single words (unigrams). Sometimes collocations can more exactly define a topic than individual words, therefore various approaches have been proposed to take into account word combinations while building topic models.

Bigram Topic Model proposed by Hanna Wallach in [11]. This model involves the introduction of the concept a hierarchical language model Dirichlet [12]. It is assumed that the appearance of a word depends on the topic and the previous word, all word pairs are collocations.

LDA Collocation Model proposed by M. Steyvers in [13]. The model introduces a new type of hidden variables x ($x = 1$, if $w_{i-1}w_i$ is collocation 0 else). This model can take into account the bigrams and unigrams, unlike the bigram topic model, where each pair of words are collocations.

N-gram Topic Model proposed by Xuerui Wang in [14]. This model adds the relation between topics and indicators of bigrams that allows us to understand the context depending on the value of the indicator [14].

PLSA-SIM proposed by Michail Nodel in [2]. The algorithm takes into account the relation between single words and bigrams (PLSA-SIM). Words and bigrams are considered as similar if they have the same component word. Before the start of the algorithm, sets of similar words and collocations are pre-calculated. The original algorithm PLSA is modified to increase the weight of similar words and phrases in case of their co-occurrence in the documents of the collection.

3 Methods to estimate the quality of topic models

To estimate the quality of topic models, several metrics were proposed.

Perplexity is a measure of inconsistency of a model towards the collection of documents. The perplexity is defined as:

$$P(D, \Phi, \Theta) = \exp \left(-\frac{1}{\text{len}(D)} \log L(D, \Phi, \Theta) \right) \quad (8)$$

Low perplexity means that the model well predicts the appearance of terms in the collection. The perplexity depends on the size of a vocabulary: usually with the increase of the collection vocabulary, the perplexity is growing.

Coherence is an automatic metric of interpretability proposed by David Newman in [15]. In [15] David Newman showed that the proposed measure of coherence has the high correlation with the expert estimates of topics interpretability.

$$PMI(w_i, w_j) = \log \frac{p(w_1, w_2)}{p(w_1)p(w_2)} \quad (9)$$

The coherence of a topic is the median coherence of word pairs representing the topic, usually it is calculated for n most probable elements in the topic. The coherence of the model is the median of the topics coherence.

A measure of the kernels uniqueness Human-constructed topics usually have unique kernels, that is words having high probabilities in the topic. The measure of kernel uniqueness shows to what extent topics are different to each other.

$$U(\Phi) = \frac{|\cup_t \text{kernel}(\Phi_t)|}{\sum_{t \in T} |\text{kernel}(\Phi_t)|} \quad (10)$$

If the uniqueness of the topic kernels is closer to one then we can easily distinguish topics from each other. If it is closer to zero then many topics are similar to each other, contain the same words in their kernels. In this paper the kernel of a topic means ten the most probable words in the topic.

4 Bigram Anchor Words Topic Modeling

The bag of words text representation does not take into account the order of words in documents, but, in fact, many words are used in phrases, which can form completely different topics.

Usually adding collocations as unique elements of the vocabulary significantly impairs the perplexity by increasing the size of the vocabulary, but the topic model interpretability is increased. The question arises: if it is possible to consider collocations in the Anchor Words algorithm without adding them to the vocabulary.

4.1 Extracting collocations (bigrams)

To extract collocations, we used the method proposed in [16]. The authors propose the following algorithm. If several words in a text mean the same entity then in this text these words should appear beside each other more often than separately. It was assumed that if a pair of words co-occurs as immediate neighbors more than half of their appearances in the same text box, it indicates that this pair of words is a collocation. For further use in topic models, we will be use 1000 most frequent bigrams extracted from the source text collection.

4.2 Representation of collocations in Anchor Words model

One of the known problems in statistical topic modeling is the high fraction of repeated words in different topics. If one wants to describe topics in a collection only with unigrams there are many degrees of freedom to determine the topics. Multiword expressions such as bigrams can facilitate more diverse description of extracted topics. Typically, the addition of bigrams as unique elements of a vocabulary increases the number of model parameters and degrades the perplexity. Further in the article, we put forward the modification of Anchor Words algorithm that can use the unigrams and bigrams as anchor words and improve the perplexity of the source Anchor topic model.

In the step 3 of the algorithm 1, each word w_i is mapped to vector \hat{Q}_i . The problem of finding the anchor words is the allocation of the "almost convex hull"[10] in the vectors \hat{Q}_i . Each topic has a single anchor word with corresponding vector from the set of \hat{Q}_i .

The space, which contains the vector \hat{Q}_i , has a thematic semantics, therefore each word may become an anchor, and thus may correspond to a some topic. To search anchor words means to find vectors corresponding to the basic hidden topics, so that the remaining topics are linear combination of basic topics.

Our main assumption is that in the space of word candidates onto anchor words positions (\hat{Q}), bigrams $w_i w_j$, are presented as a sum of vectors $w_i + w_j$. We prepare a set of bigrams and add vectors according to this bigrams in a set of anchor word candidates.

The search of the anchor words happen directly using the distance of each word on the current convex hull (Algorithm 2). Bigrams that have on their composition two vectors close to the borders of current convex hull are given the priority in the process of selection of anchor words. It is caused by the increase of the norm of the resultant vector in the direction of convex hull expansion. Therefore, while searching anchor words, we take into account bigrams and increase the probability of choosing a bigram as an anchor word that can be interpreted as a regularization.

The expansion of the convex hull helps to describe more words through the fixed basis. It is important to note that unreasonable extension of the convex hull can break the good properties of the model, such as interpretability. An algorithm for constructing the bigram anchor words model is shown at Algorithm 3. It differs from the original algorithm only in lines 4 and 5.

4.3 Experiments

The experiments were carried out on three collections:

1. **Banks Articles** – a collection of banking articles, 2500 documents (2000 for the train, 500 for the control), 18378 words.
2. **20 Newsgroups** – a collection of short news stories, 18846 documents (11314 for the train, 7532 for the control), 19570 words.

Algorithm 3 High Level Bigram Anchor Words**Input:** collection D , number of topics $\text{rem } |T|$, **set of bigrams C** **Output:** matrix Φ ;

-
- 1: $Q = \text{Word_Co-occurences}(D)$
 - 2: $\hat{Q} = \text{Rows_normalized}(Q)$
 - 3: $\hat{Q} = \text{Random_projection}(\hat{Q})$
 - 4: $\hat{B} = \hat{Q}_1, \dots, \hat{Q}_n, \hat{Q}_{C_{11}} + \hat{Q}_{C_{12}}, \dots, \hat{Q}_{C_{n1}} + \hat{Q}_{C_{n2}}$
 - 5: $S = \text{FindAnchorWords}(\hat{B}, |T|)$
 - 6: $\Phi = \text{RecoverWordTopic}(\hat{Q}, S)$
 - 7: **return** Φ
-

3. **NIPS** – a collection of abstracts from the Conference on Neural Information Processing Systems (NIPS), 1738 documents (1242 for the train, 496 for the control), 21358 words.

All collections have been preprocessed. The characters were brought to lowercase, characters which do not belong to Cyrillic and Latin alphabet were removed, words have been normalized (or stemmed for English collections), stop words and words with the length less than four letters were removed. Also words occurring less than 5 times were rejected. Collocations have been extracted with the algorithm described in Section 4.1. The preprocessed collections are available on the page, github.com/ars-ashuha/tmtk. In all experiments, the number of topics was fixed $|T| = 100$.

The metrics were calculated as follows:

- To calculate perplexity, the collection was divided into train and control parts. When calculating perplexity on test samples, each document was subdivided into two parts. In the first part, the vector of topics for the document was estimated, on the second part, perplexity was calculated .
- When calculating the coherence, the conditional probabilities are calculated with a window of 10 words.
- When calculating the unique kernel, ten most probable words in a topic were considered as its kernel.

The experiments were performed on the following models: PLSA (PL), Anchor Words (AW), Bigram Anchor Words (BiAW), Anchor Words and PLSA combination (AW + PL), Bigram Anchor Words and PLSA combination (BiAW + PL). The combination was constructed as follows: the topics obtained by the Anchor Word or Bigram Anchor Word algorithm, were used as an initial approximation for PLSA algorithm. In experiments, perplexity was measured on the control sample (P_{test}), coherence is denoted as (PMI), the uniqueness of the nuclei is denoted as (U). The results are shown in Table 1.

As in the experiments of the authors of the Anchor Words model, the perplexity grows (in two collections out of three), which is a negative phenomenon, but the uniqueness and interpretability of the topics also grows. The combina-

Table 1: Results of Numerical experiments

Collection	<i>Banks Articles</i>			<i>20 Newsgroups</i>			<i>NIPS</i>		
Metric	P_{test}	PMI	U	P_{test}	PMI	U	P_{test}	PMI	U
PL	2116	0.60	0.40	2155	0.31	0.40	1635	0.21	0.32
AW	2330	0.63	0.53	2268	0.38	0.41	1505	0.41	0.38
BiAW	2248	0.79	0.60	2183	0.68	0.54	1500	0.50	0.41
AW+PL	2052	0.78	0.58	2053	0.54	0.55	1434	0.52	0.46
BiAW+PL	1848	0.87	0.63	2027	0.78	0.64	1413	0.58	0.49

tion of *Anchor Words* and *PLSA* models shows the results better than *Anchor Words* or *PLSA* separately.

The Bigram Anchor Wordsmodel shows better results than the original *Anchor Words*: has lower perplexity, greater interpretability and uniqueness of kernels, but is still inferior to the *PLSA* model in perplexity. The combination of *Bigram Anchor Words* and *PLSA* models shows better results than other models; this combination has higher interpretability and uniqueness of the kernels.

It can be concluded that the initial approximation, given by the *Bigram Anchor Words* model, is more optimal in terms of achieving final perplexity and other metrics of quality. This approximation improves the sensitivity of *PLSA* to the initial approximation, which, in turn, can be formed taking into account the linguistic knowledge. Tables 2 and 3 contain examples of topics for Bank and NIPS collections.

Table 2: Examples of topics for the Bank collection

PLSA		ANW		ANW + PLSA + BI	
<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 1</i>	<i>Topic 2</i>
рынок	кредит	акция	кредит	рынок	кредит
российский	кредитный	рынок	кредитован	инвестор	замщик
инвестор	замщик	размещение	потребитель	акция	ипотечный
фондовый	кредитован	акционер	ипотечный	фондовый	кредитован
облигация	ставка	инвестор	замщик	инструмент	залог
бумага	банка	капитал	банк	биржа	портфель
инструмент	процентный	фондовый	население	облигация	задолжен
фонд	срок	биржа	клиент	сегмент	потребитель

Anchor words for unigram anchor model: *москва, налоговый, история, акция, сила, платеж, ассоциация*

Anchor words for bigram anchor models: *компания, миллион рубль, страна ес, управление, юридический лицо, российский федерация*

Table 3: Examples of topics for the NIPS collection

PLSA		ANW		ANW + PLSA + BI	
<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 1</i>	<i>Topic 2</i>	<i>Topic 1</i>	<i>Topic 2</i>
neuron	tree	neuron	tree	synaps	tree
spike	featur	synaps	decis	synapt	decis
fire	branch	synapt	branch	neuron	branch
time	thi	input	structur	hebbian	set
synapt	class	pattern	leaf	postsynapt	probabl
synaps	imag	neural	prune	pattern	prune
rate	object	activ	set	function	algorithm
input	decis	connect	probabl	activ	leaf

Anchor words for unigram anchor model: *face, charact, fire, loss, motion, cluster, tree, circuit, trajectori, word, extra, action, mixtur*

Anchor words for bigram anchor model: *likelihood, network, loss, face, ocular domain, reinforc learn, optic flow, boltzmann machin, markov*

5 Conclusion

We propose a modification of the Anchor Words topic modeling algorithm that takes into account collocations. The experiments have confirmed that this approach leads to the increase of the interpretability without deteriorating perplexity.

Accounting of collocations is only the first step to add linguistic information into a topic model. Further work will focus on the study of the possibilities of using the sentence structure of a text, as well as the morphological structure of words in the construction of topic models.

Acknowledgments. This work was supported by grant RFFI 14-07-00383A «Research of methods of integration of linguistic knowledge into statistical topic models».

References

1. Gao, W., Li, P., Darwish, K.: Joint topic modeling for event summarization across news and social media streams. In: Proceedings of the 21st ACM international conference on Information and knowledge management, ACM (2012) 1173–1182
2. Nokel, M., N, L.: The method of accounting bigram structure in topical models. Computational Methods and Programming **16** (2015) 215
3. Cheng, X., Yan, X., Lan, Y., Guo, J.: Btm: Topic modeling over short texts. Knowledge and Data Engineering, IEEE Transactions on **26**(12) (2014) 2928–2941

4. Krestel, R., Fankhauser, P., Nejdl, W.: Latent dirichlet allocation for tag recommendation. In: Proceedings of the third ACM conference on Recommender systems, ACM (2009) 61–68
5. Mei, Q., Cai, D., Zhang, D., Zhai, C.: Topic modeling with network regularization. In: Proceedings of the 17th international conference on World Wide Web, ACM (2008) 101–110
6. Vorontsov, K.: Additive regularization for topic models of text collections. In: Doklady Mathematics, Pleiades Publishing (2014) 301–304
7. Hofmann, T.: Probabilistic latent semantic indexing. In: Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval, ACM (1999) 50–57
8. Blei, D.M., Ng, A.Y., Jordan, M.I.: Latent dirichlet allocation. the Journal of machine Learning research **3** (2003) 993–1022
9. Arora, S., Ge, R., Moitra, A.: Learning topic models – going beyond svd. In: Foundations of Computer Science (FOCS), 2012 IEEE 53rd Annual Symposium on, IEEE (2012) 1–10
10. Arora, S., Ge, R., Halpern, Y., Mimno, D., Moitra, A., Sontag, D., Wu, Y., Zhu, M.: A practical algorithm for topic modeling with provable guarantees. arXiv preprint arXiv:1212.4777 (2012)
11. Wallach, H.M.: Topic modeling: beyond bag-of-words. In: Proceedings of the 23rd international conference on Machine learning, ACM (2006) 977–984
12. MacKay, D.J., Peto, L.C.B.: A hierarchical dirichlet language model. Natural language engineering **1**(03) (1995) 289–308
13. Steyvers M, G.T.: Matlab topic modeling toolbox 1.3. (2005)
14. Wang, X., McCallum, A., Wei, X.: Topical n-grams: Phrase and topic discovery, with an application to information retrieval. In: Data Mining, 2007. ICDM 2007. Seventh IEEE International Conference on, IEEE (2007) 697–702
15. Newman, D., Lau, J.H., Grieser, K., Baldwin, T.: Automatic evaluation of topic coherence. In: Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, Association for Computational Linguistics (2010) 100–108
16. B, D., N, L., S, S.: Forming the base of terminological phrases in the texts of the subject area. Tp (2003) 201–210